# Lecture 20: Belief propagation:

Def : (Undirected) graphical model w/ pairwise interactions:

Let $\{\psi_{ij}\}_{(i,j) \in F}$ be <u>compatibility functions</u> $\{\pm 1\}^2 \to \mathbb{R}_{\geq 0}$ that dictate interactions b/t particles

The <u>Gibbs measure</u> is dist. over $\{\pm 1\}^n$ given by

$$\mu(\underset{\text{"spins"}}{x}) \triangleq \frac{1}{Z} \prod_{(ij) \in F} \psi_{ij}(x_i, x_j), \qquad Z \triangleq \sum_{x \in \{\pm 1\}^n} \prod_{(ij)} \psi_{ij}(x_i, x_j)$$

where $Z$ is the <u>partition function</u>, i.e. normalizing constant.

We will use the shorthand
$$\mu \propto \prod_{(ij)} \psi_{ij}(x_i, x_j)$$

---

## Example ("Ising model"):

$$\psi_{ij}(x_i, x_j) = \exp(-\beta A_{ij} x_i x_j), \text{ so}$$

$$\mu(x) \propto \exp\left(-\frac{\beta}{2} \underbrace{x^T A x}_{\text{"energy"}}\right)$$

for $A \in \mathbb{R}^{n \times n}$ a symmetric matrix with zero diagonal

$\beta$: "inverse temperature"

$A$: "Hamiltonian" / "interaction matrix"

$$\mathcal{E}(x) \triangleq \frac{\beta}{2} x^T A x = -\lg\left(\prod_{(ij) \in F} \psi_{ij}(x_i, x_j)\right) : \text{"energy"}$$
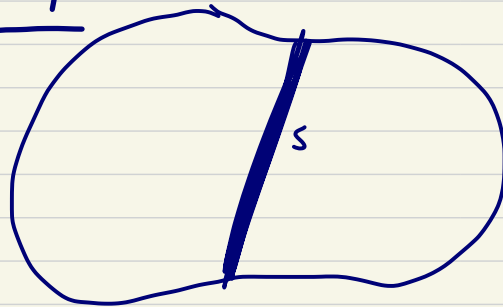
As $\beta \to 0$, $\mu \to \text{Unif}(\{\pm 1\}^n)$
$\beta \to \infty$, $\mu \to \text{Unif}(\{\text{energy minimizers}\})$

Can think of $-A$ as adjacency matrix of weighted graph.* Denote this by $G$.

$$\partial i \triangleq \{j \text{ s.t. } A_{ij} \neq 0\} = \{j \text{ s.t. } (i,j) \in F\},$$

i.e. the neighbors of node $i$ in $G$

## Markov property:



If $[n] \setminus S$ decomposes into disjoint pieces, then marginal distributions on the pieces are <u>independent</u>, conditioned on any assignment to the spins on $S$

e.g. if we condition on $x_{\partial i}$, then conditional dist. on $x_i$ is independent of conditional dist on rest of the spins. For Ising model:

$$\Pr\left[x_i = \sigma \mid x_{\partial i} = s\right] \propto \exp\left(2\beta \sum_{j \in \partial i} A_{ij} s_j \sigma\right)$$

2 fundamental algorithmic tasks in inference:

① computing the partition function $Z$

② Sampling from Gibbs measure $\mu$

Note  alg. for ① $\Rightarrow$ alg. for ② and vice versa
("equivalence of counting + sampling")

<u>Challenge</u>: $Z$ is sum of exponentially many terms, so
in many cases we expect it is computationally hard to
compute...

e.g. if $\psi_{ij}(x_i, x_j) = \mathbb{1}[x_i \neq x_j]$ for all $(i,j) \in E(G)$,

$Z = \#$ independent sets of $G$  ("#P-complete",
                                                              i.e. <u>very</u> hard)

So our goal will be to approximate $Z$ / approximately
from Gibbs measure $\mu$

Some approaches :

- Markov chain Monte Carlo (MCMC)
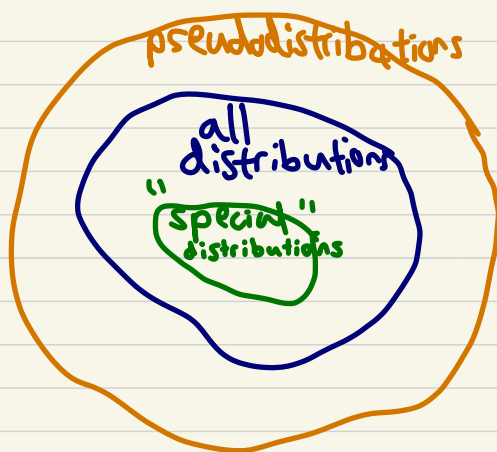
✗ - Variational inference (VI)

- Diffusion models ( very recent, more on this at the
                                        end of the course)

<u>VI</u> : approximate $\mu$ by dist. from family $p$
of "simpler" distributions that are easy
to sample from (e.g. product distributions,
aka mean-field!) :

$$\min_{\nu \in p} KL(\nu \| \mu) \qquad (\clubsuit)$$

Note: - if $p$ is <u>all</u> distributions, minimizer is $\nu^{\clubsuit} = \mu$
(Gibbs' inequality)

- "opposite" of SoS relaxation



An issue: can't even evaluate the objective function
in $(\clubsuit)$, let alone optimize!

Fortunately, this particular issue is not really an issue:

$$KL(\nu \| \mu) = \mathbb{E}_\nu \ln \frac{\nu}{\mu}$$

$$= \mathbb{E}_\nu \ln \frac{\nu}{e^{-\mathcal{E}}/z}$$

free energy of the Gibbs measure

$$= \mathbb{E}_\nu \ln \nu + \mathbb{E}_\nu \mathcal{E} - \ln(1/z)$$

$\underbrace{\quad}$ negative entropy $-H(\nu)$

$\underbrace{\quad}$ if $\nu$ simple, easy to evaluate

$\underbrace{\quad}$ average energy

$\underbrace{\quad}$ easy to approximate

independent of $\nu$!

$$G[\nu] \doteq -H(\nu) + \mathbb{E}_\nu[\mathcal{E}]$$

"Gibbs free energy functional" /

$-1 \times$ "evidence lower bound" (ELBO)

so to minimize $KL(\nu \| \mu)$, suffices to minimize $G[\nu]$ which is easy to evaluate

Interpretation of $G$ as "regularized energy": for Ising model, recall $\mathcal{E}(x) = \frac{\beta}{2} x^T A x$, so

$$G[\nu] = \frac{\beta}{2} x^T A x \underbrace{- H(\nu)}_{\text{"entropy regularization"}}$$

"hot"

When $\beta$ small, minimizer prioritizes maximizing entropy

$\beta$ big, minimizer prioritizes minimizing avg. energy

"cold"

$G[\nu]$ efficiently computable, but computationally intractable to optimize a priori...

Rest of lecture: powerful heuristic, <u>belief propagation</u> (BP), for solving $\min_\nu G[\nu]$.

2 interpretations of the heuristic:

① dynamic programming

② finding stationary points of a relaxation of the Gibbs free energy (Bethe free energy) (see supplemental notes)

## <u>BP as dynamic programming</u>:

Let's first shift focus to easier task than full-blown VI: <u>marginal estimation</u> dist. $\mu_i$ over each node is a Bernoulli random variable, goal is to estimate it
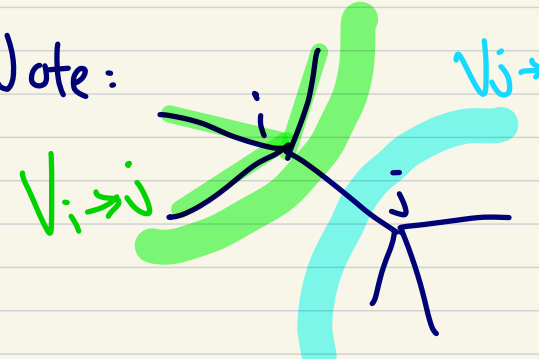
<u>Physics motivation</u>:
physicists care about limiting objects, and one important that they consider as $n \to \infty$ is the empirical dist. over marginals, i.e.

$$q_n(z) \triangleq \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}\left[z = \mu_i\right]$$

and, given a sequence of Gibbs measures $(\mu^{(n)})$, want to understand $\lim_{n \to \infty} q_n$

To motivate the algorithm, assume $G$ is a tree

Note: 

$V_{j \to i}$ removing $(i,j)$ from tree splits $G$ into two subtrees

$V_{i \to j}$

$\bar{V}_{i \to j} \quad (V_{j \to i} + \text{edge } (i,j))$

$\bar{V}_{j \to i} \quad (V_{i \to j} + \text{edge } (i,j))$

To sample from $\mu_i$,

1). Sample spins on subtrees $V_{j \to i}$ for $j \in \partial i$, yields assignment $s \in \{\pm 1\}^{\partial i}$ to $\partial i$

2). Sample from conditional dist. on $x_i$, i.e.

$$\Pr\left[x_i = \sigma \mid x_{\partial_i} = s\right] \propto \prod_{j \in \partial i} \psi_{ij}(\sigma, s_j)$$

By law of total probability,

$(✿✿)$ 
$$\Pr[x_i = \sigma] \propto \sum_{S \in \{\pm 1\}^{\partial i}} \prod_{j \in \partial i} \Pr[x_j = s_j] \, \psi_{ij}(\sigma, s_j)$$

$$= \prod_{j \in \partial i} \underbrace{\sum_{s_j \in \{\pm 1\}} \Pr[x_j = s_j] \, \psi_{ij}(\sigma, s_j)}_{M_{V_{j \to i}}}$$

b/c marginal dist's on spins in $V_{j \to i}$ are independent across $j$'s

$(✿✿✿)$

proportional to $\Pr[x_i = \sigma]$
$$\underset{M_{\bar{V}_{j \to i}}}{}$$

(i.e. can express marginals of $M_{\bar{V}_{j \to i}}$ in terms of marginals of $M_{V_{j \to i}}$)

Unsatisfying b/c we've gone from

$\Pr_{M}[x_i = \sigma]$ to $\Pr_{M_{\bar{V}_{j \to i}}}[x_i = \sigma]$, but

we're very close.

Define messages :

$$m_\sigma^{(j)\to i} \;\triangleq\; \Pr_{\mu_{V_{j\to i}}}\left( x_j = \sigma \right)$$

$$\overline{m}_\sigma^{\,j\to(i)} \;\triangleq\; \Pr_{\mu_{\overline{V}_{j\to i}}}\left( x_i = \sigma \right)$$

Then (★☆☆) can be written as

$$\boxed{\;\overline{m}_\sigma^{\,j\to(i)} \;\propto\; \sum_{s\in\{\pm 1\}} m_s^{(j)\to i} \cdot \psi_{ij}(\sigma, s)\;} \qquad \text{I}$$

Also note that (★☆☆) can be modified to apply to $\mu_{V_{i\to k}}$ instead of $\mu$, i.e.

previously, $(\cancel{666})$ gave

$$\Pr_{M}\left[x_i = \sigma\right] \propto \prod_{j \in \partial i} \bar{m}_\sigma^{j \to i}$$

after removing edge $(i,k)$, we get

(II)
$$\boxed{m_\sigma^{i \to k} \propto \prod_{j \in \partial i \setminus \textcolor{red}{k}} \bar{m}_\sigma^{j \to i}}$$

$$\left( = \Pr_{M_{V^{i \to k}}}\left[x_i = \sigma\right] \right)$$

$\searrow$ we can then write marginals succinctly
in terms of the messages:

$(\cancel{66666})\quad \Pr_{M}\left(x_i = \sigma\right) = m_\sigma^{i \to j}\, m_\sigma^{j \to i}$

Combining (I) and (II) yields:

(rec)
$$\boxed{m_\sigma^{i \to k} \propto \prod_{j \in \partial i \setminus k} \sum_{s \in \{\pm 1\}} m_s^{j \to i} \cdot \psi_{ij}(\sigma, s)}$$

BP on trees:

1). Pick arbitrary root vertex
2). For every leaf $j$ and parent $i$, initialize $m_\sigma^{j \to i} = 1/2 \quad \forall \sigma \in \{\pm 1\}$
3). Use (rec) to compute $\bar{m}$'s via dynamic programming, starting from leaves
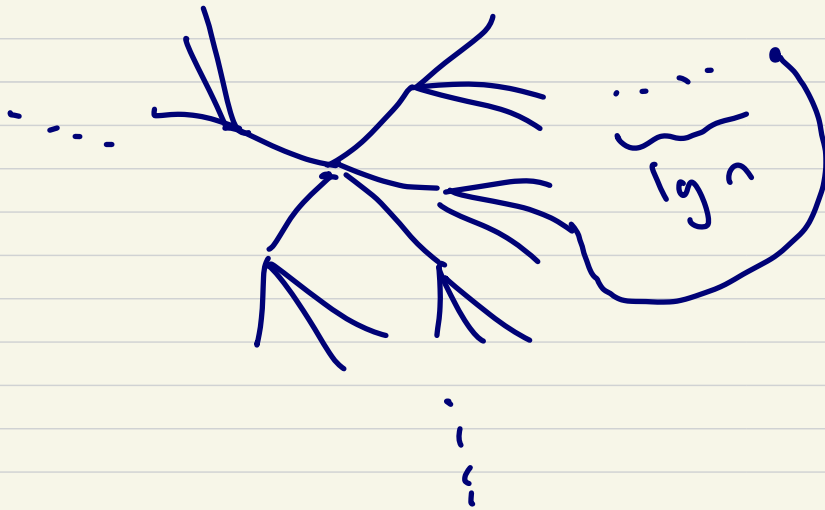4) Use ① to compute $m$'s
5) Use (✦✦✦) to compute marginals

What if $G$ is not a tree? Then subtree marginals $\{\mu_v^{j \to i}\}_{j \in \partial i}$ are not independent...

Nevertheless, can still run the above algorithm* and hope it gives something interesting!

* As stated, the algorithm is stated w/ a tree structure in mind. Without this, we can still apply update rules for $\bar{m}$ and $m$ in parallel

Intuition for why this is a good idea:
if the graph is a random sparse graph,
then <u>locally</u> it looks like a tree



If every edge appears w.p. $\frac{c}{n}$ for $c = O(1)$,
then probability that some "descendant" at depth $d$
"returns" to ancestor $i$ is

$$1 - \left(1 - \frac{1}{n}\right)^{c^d}$$

so as long as $c^d \ll n$, this is $o(1)$.

Next lecture we will see a natural setting where
such a sparse random graph arises.

Even in such cases, BP is notoriously hard to analyze. We will instead see 2 rigorous alg's inspired by BP:

1). Spectral methods on nonbacktracking operators

2). approximate message passing.